

# AI & ML

---

SESSION 20  
MODULE 5

# Session Outline

---

- Curse of Dimensionality
- Dimensionality Reduction
- Feature Selection
- Matrix Factorization based DR
- Preliminaries
- PCA

# Can you rank the models ?

---

Model 1  
D = 10  
Accuracy a1

Model 2  
D = 50  
Accuracy a2

Model 3  
D = 100  
Accuracy a3

Model 5  
D = 200  
Accuracy a4

Model 5  
D = 500  
Accuracy a5

# Can you rank the models ?

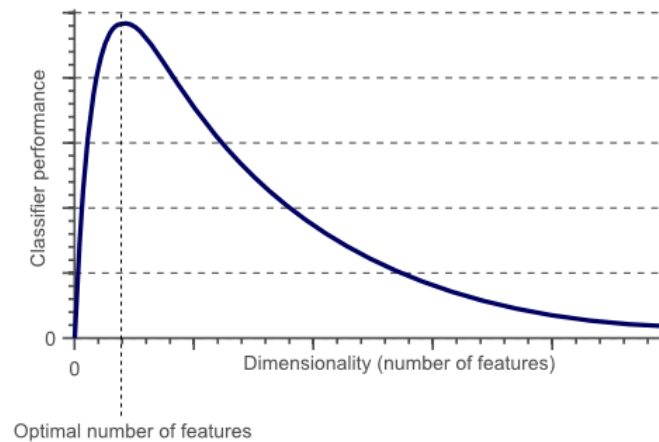
Model 1  
D = 10  
Accuracy a1

Model 2  
D = 50  
Accuracy a2

Model 3  
D = 100  
Accuracy a3

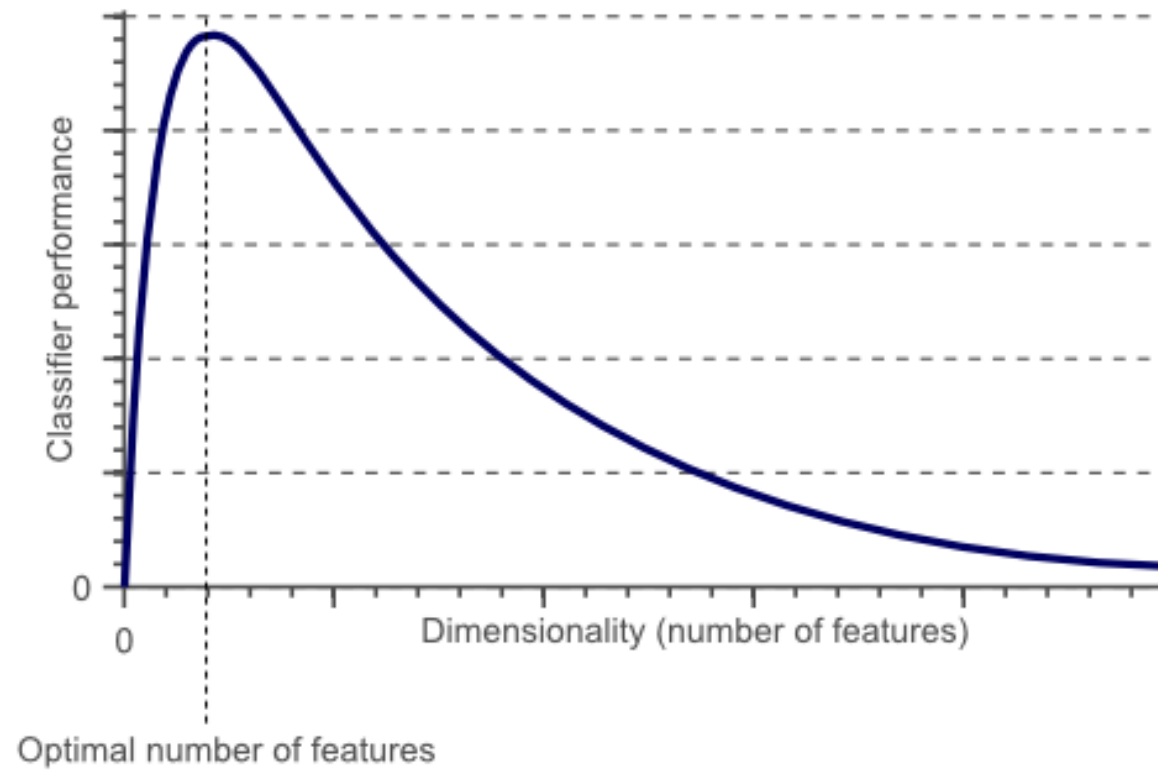
Model 5  
D = 200  
Accuracy a4

Model 5  
D = 500  
Accuracy a5



# Hughes Phenomenon

---



# Curse of Dimensionality

---

- Term coined by Bellman in 1961.
- The sample size needed to estimate a function of several variables to a given degree of accuracy (i.e. to get a reasonably low-variance estimate) grows exponentially with the number of variables.
- as the number of attributes or the dimensions increases, the number of training samples required to generalize a model also increases phenomenally.
- high-dimensional spaces are inherently sparse.
- Overfitting models.

# Combating CoD

---

A way to avoid the curse of the dimensionality is to reduce the input dimension of the function to be estimated.

# Dimensionality Reduction

---

- ❖ Dimensionality – no of input variables.
- ❖ Large numbers of input features can cause poor performance for machine learning algorithms.
- ❖ Dimensionality reduction is a general field of study concerned with reducing the number of input features.
- ❖ Dimensionality reduction methods include feature selection, linear algebra methods, projection methods, and autoencoders.



# When to do

---

- ☐ Performed after data cleaning and data scaling
- ☐ Before training a predictive model.

Any dimensionality reduction performed on training data must also be performed on new data, such as a test dataset, validation dataset, and data when making a prediction with the final model.

# Common Data Preparation Tasks

---

- ☐ Data Cleaning
- ☐ Feature Selection
- ☐ Data Transforms
- ☐ Feature Engineering
- ☐ Dimensionality Reduction

# Common Data Preparation Tasks

---

☐ Data Cleaning

☐ **Feature Selection**

☐ Data Transforms

☐ Feature Engineering

☐ **Dimensionality Reduction**

- ❖ Reduce computational cost of modeling
- ❖ Improve model performance

# Techniques for Dimensionality Reduction

---

# Feature Selection Methods

---

# Feature Selection

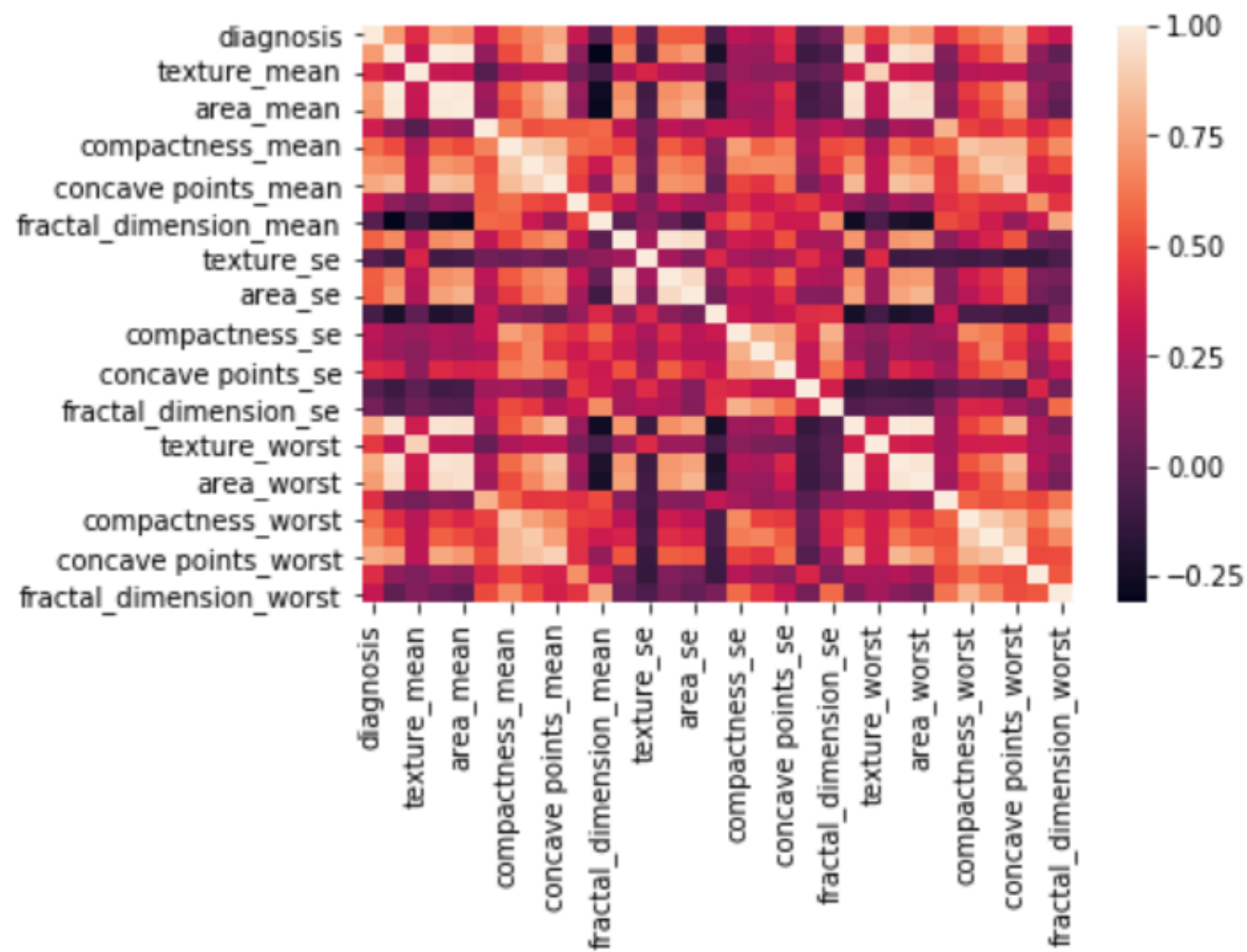
---

Process of reducing the number of input variables when developing a predictive model.

Already seen

- Delete columns with id, names.
- Delete columns with constant values.
- Delete columns that exhibit low variance.

# Correlation



# Score based Selection

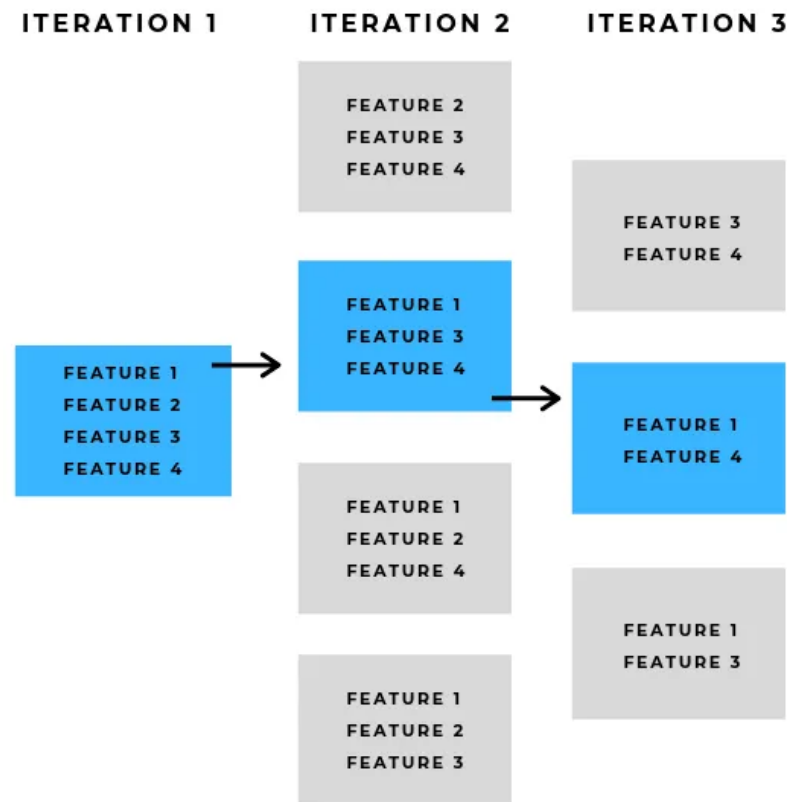
---

Selecting K best features

[sklearn.feature\\_selection.SelectKBest — scikit-learn 1.2.1 documentation](https://scikit-learn.org/stable/modules/generated/sklearn.feature_selection.SelectKBest.html)

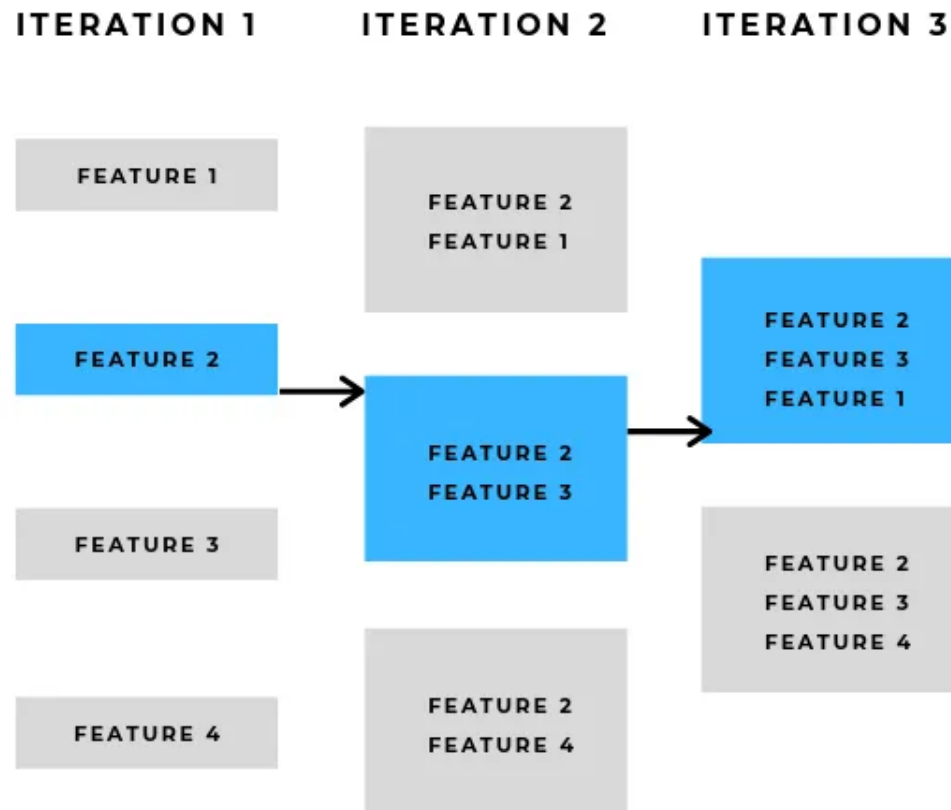


# Recursive Feature Elimination



<https://medium.com/analytics-vidhya/feature-selection-methods-for-data-science-just-a-few-fca3086eb445>

# Sequential Feature Selection

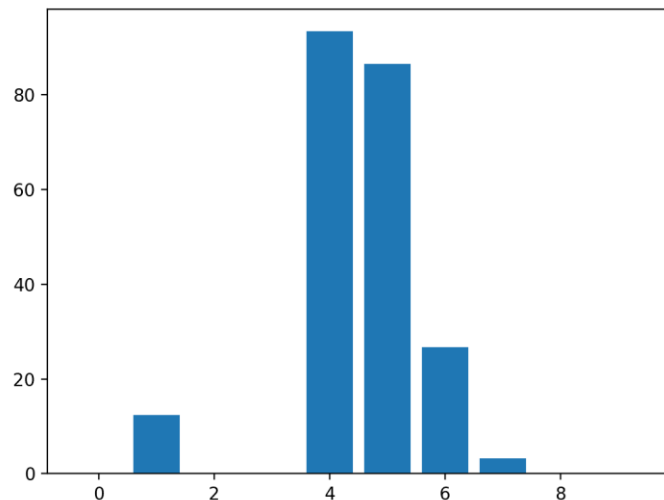


<https://medium.com/analytics-vidhya/feature-selection-methods-for-data-science-just-a-few-fca3086eb445>

# Feature Importance based selection

Linear machine learning algorithms fit a model where the prediction is the weighted sum of the input values.

All of these algorithms find a set of coefficients to use in the weighted sum in order to make a prediction. These coefficients can be used directly as a crude type of feature importance score.



Linear Regression Coefficients as Feature Importance Scores

Feature: 0, Score: 0.00000  
Feature: 1, Score: 12.44483  
Feature: 2, Score: -0.00000  
Feature: 3, Score: -0.00000  
Feature: 4, Score: 93.32225  
Feature: 5, Score: 86.50811  
Feature: 6, Score: 26.74607  
Feature: 7, Score: 3.28535  
Feature: 8, Score: -0.00000  
Feature: 9, Score: 0.00000

# Matrix Factorization

---

# Matrix Factorization

---

Matrix factorization methods can be used to reduce a dataset matrix into its constituent parts.

Examples include the eigen decomposition and singular value decomposition.

The parts can then be ranked and a subset of those parts can be selected that best captures the salient structure of the matrix that can be used to represent the dataset.

The most common method for ranking the components is principal components analysis, or PCA for short.

# Preliminaries

---

# Variance

---

A measure of the spread of the data in a data set with mean,  $\bar{x}$

$$\sigma^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{(n-1)}$$

Variance is claimed to be the original statistical measure of spread of data.

# Covariance

---

Variance – measure of the deviation from the mean for points in one dimension, e.g., heights

Covariance – a measure of how much each of the dimensions varies from the mean with respect to each other.

Covariance is measured between 2 dimensions to see if there is a relationship between the 2 dimensions, e.g., number of hours studied and grade obtained.

The covariance between one dimension and itself is the variance

$$\text{var}(X) = \frac{\sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})}{(n-1)}$$
$$\text{cov}(X, Y) = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{(n-1)}$$



# Covariance Matrix

---

Representing covariance among dimensions as a matrix, e.g., for 3 dimensions:

$$C = \begin{bmatrix} \text{cov}(X, X) & \text{cov}(X, Y) & \text{cov}(X, Z) \\ \text{cov}(Y, X) & \text{cov}(Y, Y) & \text{cov}(Y, Z) \\ \text{cov}(Z, X) & \text{cov}(Z, Y) & \text{cov}(Z, Z) \end{bmatrix}$$

Properties:

- ❑ Diagonal: variances of the variables
- ❑  $\text{cov}(X, Y) = \text{cov}(Y, X)$ , hence matrix is symmetrical about the diagonal (upper triangular)
- ❑ m-dimensional data will result in mxm covariance matrix

# Linear Independence

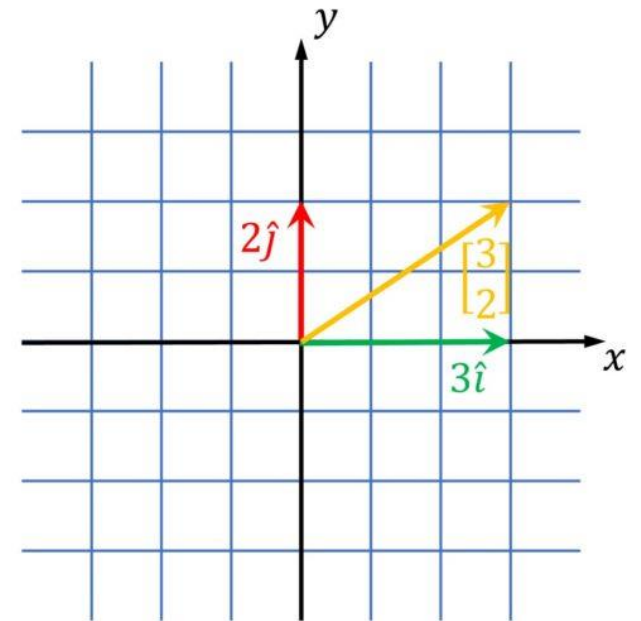
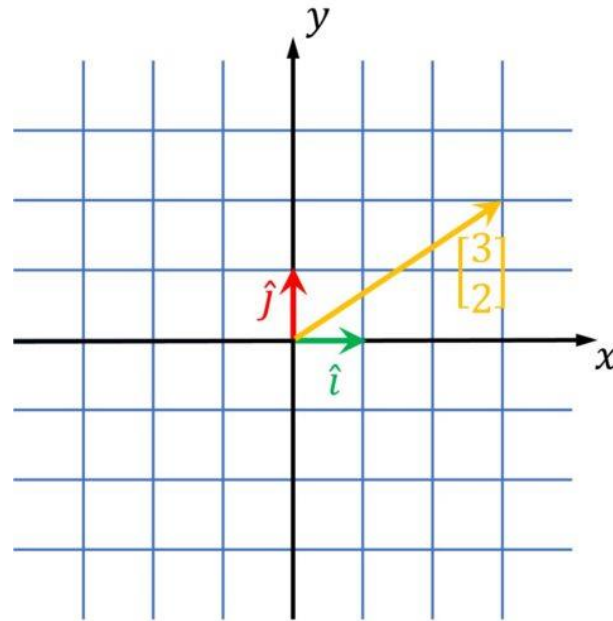
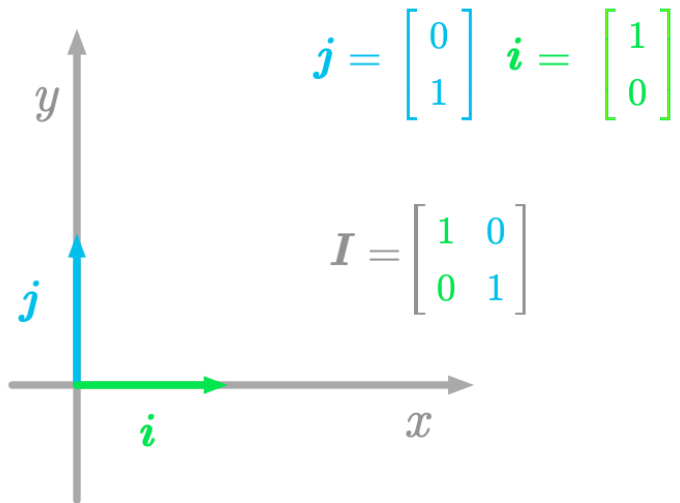
---

A set of  $n$ -dimensional vectors  $x_i \in \mathbb{R}^n$ , are said to be linearly independent if none of them can be written as a linear combination of the others.

$$c_1x_1 + c_2x_2 + \dots + c_kx_k = 0$$

*iff*  $c_1 = c_2 = \dots = c_k = 0$

# Basis



?

---

$$\begin{bmatrix} 2 & 3 \\ 2 & 1 \end{bmatrix} \times \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

# Example

---

$$\begin{bmatrix} 2 & 3 \\ 2 & 1 \end{bmatrix} \times \begin{bmatrix} 3 \\ 2 \end{bmatrix} = \begin{bmatrix} 12 \\ 8 \end{bmatrix} = 4 \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

$$A \times v = \lambda v$$

# Eigen values and Eigen vectors

---

$$A \times v = \lambda v$$

A: m x m matrix

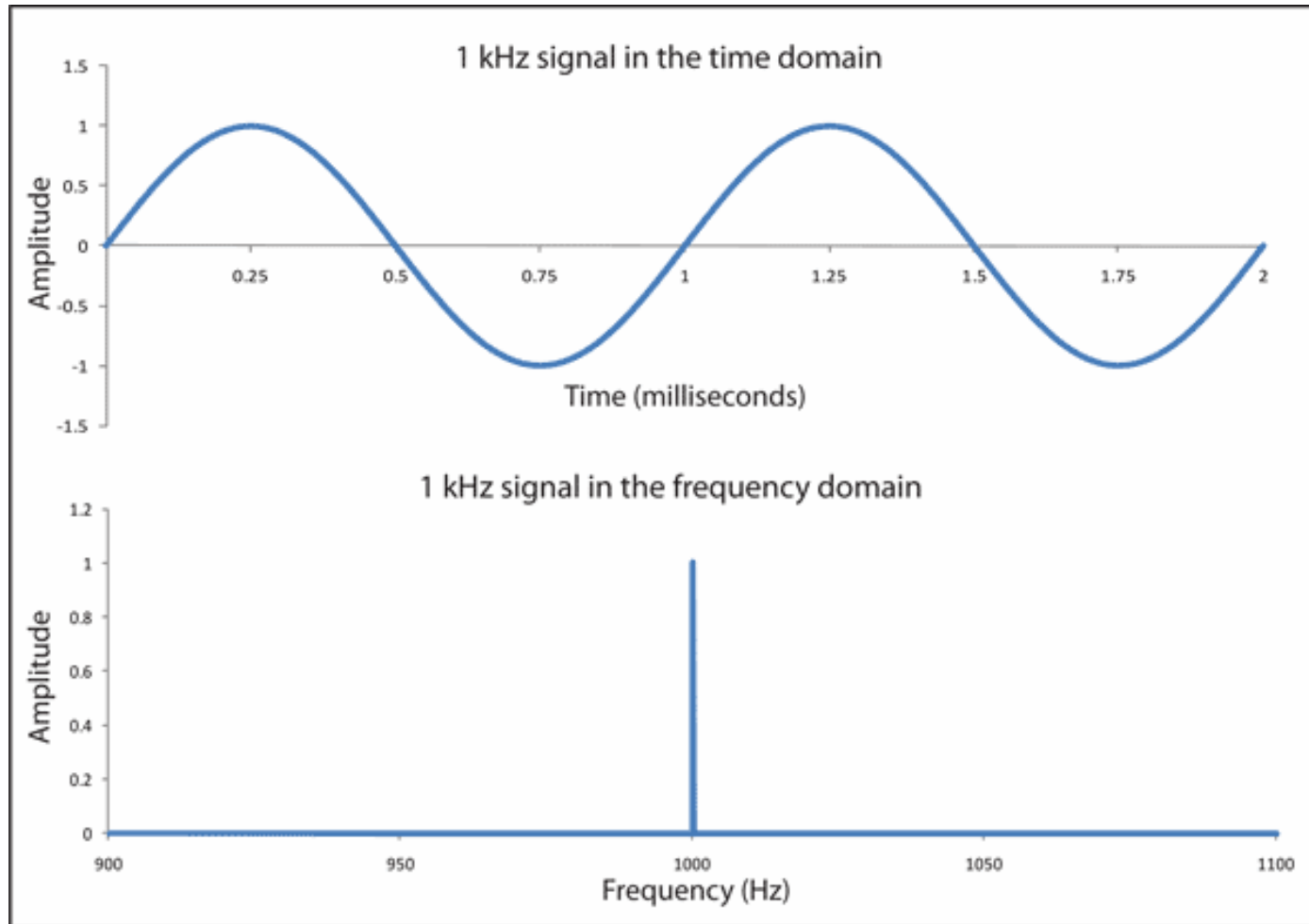
v: m x 1 non-zero vector

$\lambda$ : scalar

Any value of  $\lambda$  for which this equation has a solution is called the eigenvalue of A and the vector v which corresponds to this value is called the eigenvector of A.

# Principal Component Analysis

---





# Change of Basis

---

Let  $X$  and  $Y$  be  $m \times n$  matrices related by a linear transformation  $P$ .

$X$  is the original recorded data set and  $Y$  is a re-representation of that data set.

$$PX = Y$$

Let's define;

- $p_i$  are the rows of  $P$ .
- $x_i$  are the columns of  $X$ .
- $y_i$  are the columns of  $Y$ .

$$PX = Y$$

---

- P is a matrix that transforms X into Y.
  - Geometrically, P is a rotation and a stretch (scaling) which again transforms X into Y.
  - The rows of P,  $\{p_1, p_2, \dots, p_m\}$  are a set of new basis vectors for expressing the columns of X.
- 
- *Changing the basis doesn't change the data – only its representation.*
  - *Changing the basis is actually projecting the data vectors on the basis vectors.*

# Change of basis problem

---

- Assuming linearity, the problem now is to find the appropriate change of basis.
- Now,
  - What is the best way to re-express  $X$ ?
  - What is the good choice of basis  $P$ ?

# Principal Component Analysis

---

a statistical procedure to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables.

Each of the principal components is chosen in such a way so that it would describe most of them still available variance .

All these principal components are orthogonal to each other.

In all principal components first principal component has a maximum variance.

These are basically performed on a square symmetric matrix.

# Steps in PCA

---

*Step 1:* Standardization of the values of all the features in the dataset to zero mean and unit variance.



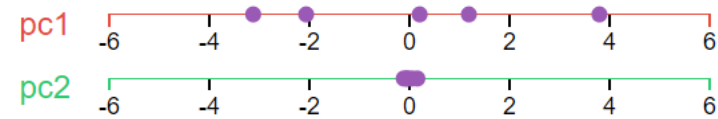
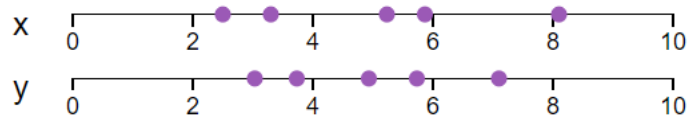
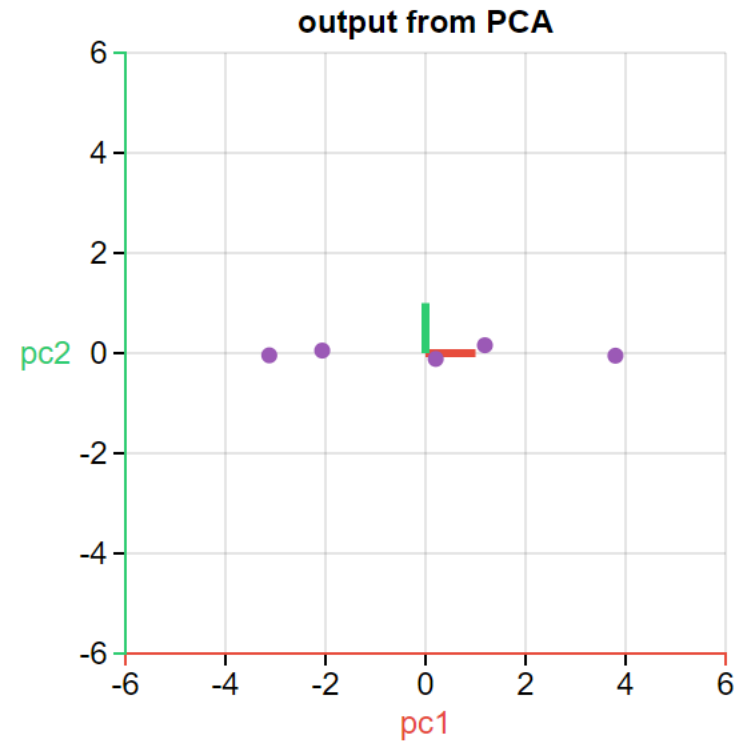
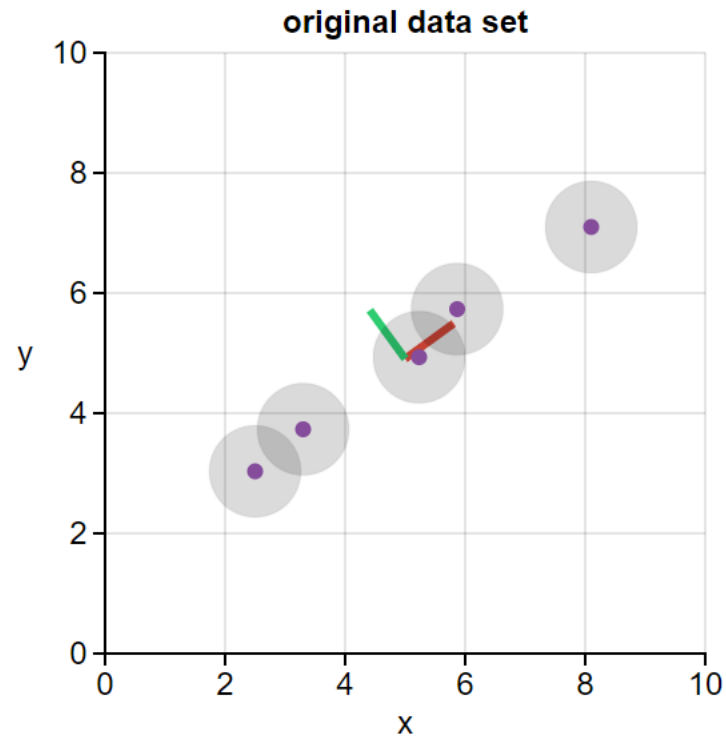
*Step 2:* Obtaining the covariance matrix from the training dataset.



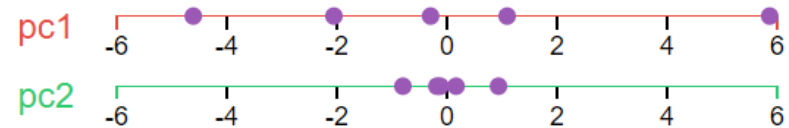
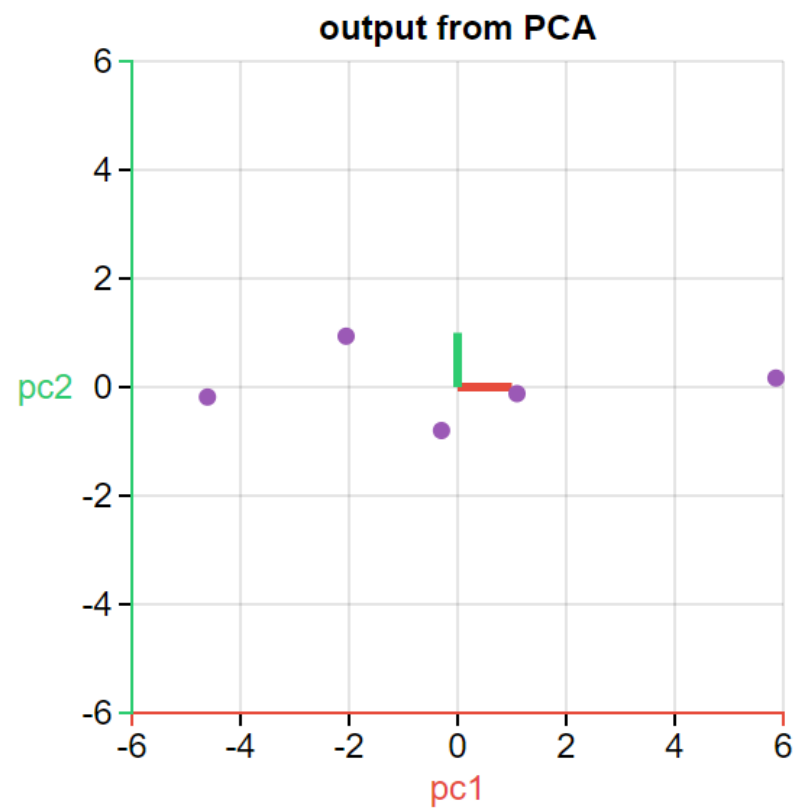
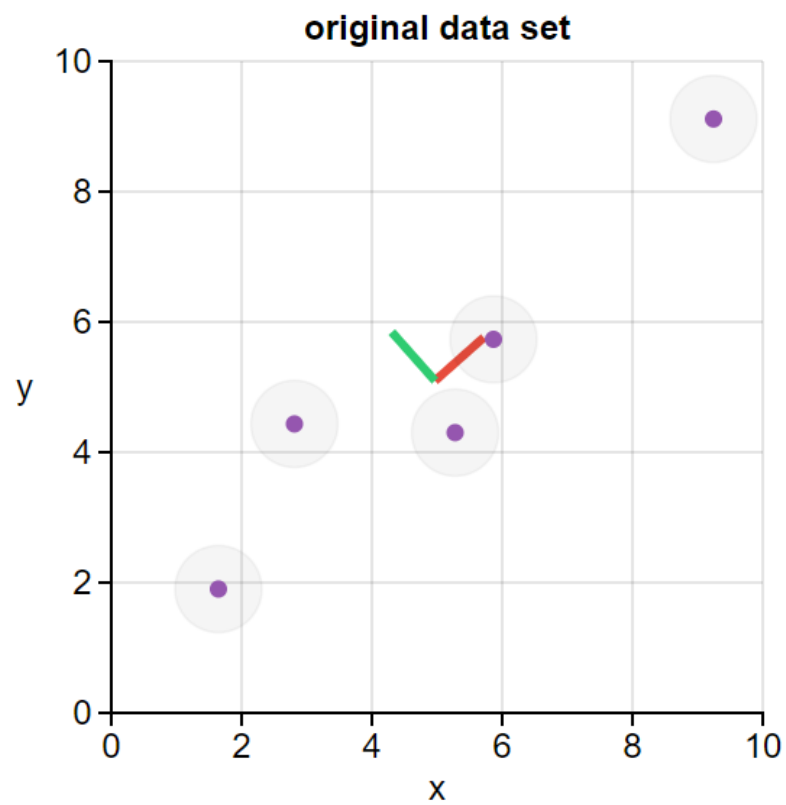
*Step 3:* Obtaining the Eigen values and Eigen vectors from the covariance matrix.



*Step 4:* Project the data points in the testing dataset in the direction of PCs of the training dataset.



[Principal Component Analysis explained visually \(setosa.io\)](https://setosa.io/Principal-Component-Analysis-explained-visually)



# PCA on a dataset

---



# Thank You!

---